

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

Recognizer of Text-Based Work

Inventors:

Ramarathnam Venkatesan

Michael Malkin

ATTORNEY'S DOCKET NO. MS1-647US

TECHNICAL FIELD

This invention relates to a technology for recognizing a text-based work.

BACKGROUND

Detecting and determining the existence of text plagiarism is complex and difficult. This complexity and difficulty increases in direct proportion with the amount of available text documents. This is the age of electronic commerce, with so-called “e-books”, Internet, HTML, e-mails, cyber-classrooms and textbooks, electronic publishing, electronic faxing, scanning, Portable Document Format (pdf) documents, Web pages, newspapers on-line, optical character recognition (OCR), “cut-and-paste”, and pay-per-chapter electronic publishing, etc. becoming common place. In this age, text, copies of text, copies of copies of text, etc. fly across the world in a matter of seconds.

In this age, it is thoughtlessly commonplace to electronically copy text and do so instantaneously with a click of a button. It is exceedingly easy to duplicate wholesale (or significant) portions of text documents. This task requires no more technical expertise than the ability to press a button or press CNTL-V (to complete the “cut-and-paste” operation).

Plagiarism

However, just because it is easy to do something, does not make it right. Although it is easy for a person to copy an author’s work and pawn it off as his own, it does not make such action right. Such action is commonly called “cheating” or “plagiarism.” Thus, a person engaging in such action is a “cheater”

1 or a “plagiarizer.” Since most contemporary works are copyrighted (either
2 automatically or upon registration), a plagiarizer is also infringing such copyrights
3 and is subject to civil and possibly criminal penalties.

4 Why would a plagiarizer take action that is socially unacceptable, deceitful,
5 and likely illegal? It is easy for the plagiarizer to do and it is unlikely for him to
6 be caught.

7 A plagiarizer realizes that authorities must compare the pilfered words in
8 his work with oceans of words, phrases, quotes, chapters, books, and other works.
9 These oceans are vast and deep. The oceans include text found in all of the
10 libraries, bookstores, web sites, manuals, textbooks, e-mails, etc. of the whole
11 world.

12 Catching a plagiarizer is a daunting task indeed. Typically, if an
13 investigative authority does not have a lead for a place to look, it is nearly
14 impossible. However, one tool that makes the investigation easier is an electronic
15 database (or index) of text that has been recorded electronically.

16 To avoid capture, a plagiarizer may simply change a few token words,
17 punctuations, pagination, text order, insertion of new text, and/or format in the text
18 documents. Meanwhile, the true authors and publishers of the substantive content
19 of the plagiarized work are robbed of well-deserved credit and/or royalties.

20 **Conventional Efforts to Detect Plagiarism**

21
22 Much effort has been directed towards protecting images, audio, and video
23 by either embedding a hidden watermark and/or generating a mathematical
24 representation of such content. Much of this effort is geared towards detecting
25 identifiers within the content even after the signals have been modified

(intentionally or purposefully). Such identifiers may be inserted into the content or be inherent in the content.

Generally, these conventional techniques may insert an imperceptible change in multimedia (such as audio or video). Alternatively, these techniques determine an inherent characteristic of a work. These conventional techniques rely on the foundation that the code/inherent characteristic cannot be detected without access to secret knowledge (such as a cryptographic key) and is unalterable without noticeably altering the content.

However, these conventional techniques have not been directed toward protecting text because they do not apply to text. They don't apply to text because these conventional techniques generally require a perceptual change to the original content or they are easily thwarted.

For example, the concept of embedding a watermark into an image or audio signal does not apply to text because embedding a watermark would significantly alter the content—unless, of course, the author inserts it. That alteration would be clearly perceivably noticeable. A mathematical representation of text is easily thwarted by changing a few token words, punctuations, insertion of new text, pagination, text order, and/or format in the text documents

Side-by-side Text Comparison Approach. A side-by-side comparison of suspect text and possibly original text is an existing technique for detecting a copy of an original text. However, it can be easily thwarted by reordering text, adding text, and changing un-essential text. If a comparison is done manually, a person may overlook such obfuscation tactics and see through to the similarity (which may amount to plagiarism). However, a comparison of electronic documents by a computer is not so forgiving.

1 With the emergence of so-called e-books, the problem of protecting text is
2 becoming more important. E-books refer to the electronic distribution of
3 electronic text. It is an alternative commercial publication technique.

4 Although such e-book mechanisms include cryptographic locks, such locks
5 can be picked. Although no conventional technique is available, it would be
6 helpful to determine if a subject body of text is substantially similar to an original
7 text.

8 9 **Content Categorization**

10 Like plagiarism, categorizing the content of a text-based work often
11 requires a subjective comparison of existing works. Works of similar nature are
12 grouped into the same category. Text-based works may be classified into any
13 number of categories, such as mystery novels, math textbooks, non-fiction books,
14 self-help books, commercial web pages, poetry, and the other such works.

15 Typically, such categorization is subjectively determined by manual (i.e.,
16 human) subjective analysis of a work so that it may be grouped with an existing
17 category. No such technique exists for automatically (i.e., without substantial
18 human involvement) analyzing and categorizing a text-based work.

19 20 **SUMMARY**

21 Described herein is a technology for recognizing the content of text
22 documents. The technology may detect similarity between text-based works in an
23 automatic and accurate manner. Furthermore, it may categorize content of text-
24 based works in an automatic and accurate manner.

1 Generally, the technology determines one or more hash values for the
2 content of a text document. Furthermore, the technology may generate a "sifted
3 text" version of a document.

4 In one implementation described herein, document recognition is used to
5 determine whether the content of one document is copied (i.e., plagiarized) from
6 another document. This is done by comparing hash values of documents (or
7 alternatively their sifted text).

8 In another implementation described herein, document recognition is used
9 to categorize the content of a document so that it may be grouped with other
10 documents in the same category.

11 This summary itself is not intended to limit the scope of this patent.
12 Moreover, the title of this patent is not intended to limit the scope of this patent.
13 For a better understanding of the present invention, please see the following
14 detailed description and appending claims, taken in conjunction with the
15 accompanying drawings. The scope of the present invention is pointed out in the
16 appending claims.

17 **BRIEF DESCRIPTION OF THE DRAWINGS**

18
19 The same numbers are used throughout the drawings to reference like
20 elements and features.

21 Fig. 1 is a schematic block diagram showing an embodiment of a
22 recognizer of text-based work.

23 Fig. 2 is a schematic block diagram showing another embodiment of a
24 recognizer of text-based work.
25

Fig. 3 is a schematic block diagram showing still another embodiment of a recognizer of text-based work.

Fig. 4 is a flow diagram showing an illustrative methodological implementation of a recognizer of text-based work.

Fig. 5 is a flow diagram showing another methodological implementation of a recognizer of text-based work.

Fig. 6 is a flow diagram showing still another methodological implementation of a recognizer of text-based work.

Fig. 7 is an example of a computing operating environment capable of implementing an implementation of a recognizer of text-based work.

DETAILED DESCRIPTION

The following description sets forth specific embodiments of the recognizer of text-based work that incorporate elements recited in the appended claims. These embodiments are described with specificity in order to meet statutory written description, enablement, and best-mode requirements. However, the description itself is not intended to limit the scope of this patent. Rather, the inventors have contemplated that the claimed present invention might also be embodied in other ways, in conjunction with other present or future technologies.

Incorporation by Reference

The following co-pending patent applications are incorporated by reference herein (which are all assigned to the Microsoft Corporation):

- U.S. Patent Application Serial No. 09/390271, entitled “A Technique for Watermarking an Image and a Resulting Watermarked Image” filed Sept. 7, 1999;
- U.S. Patent Application Serial No. 09/390272, entitled “A Technique for Detecting a Watermark in a Marked Image” filed on Sept. 7, 1999;
- U.S. Patent Application Serial No. _____, entitled “Robust Recognizer of Perceptually Similar Content” filed on April 24, 2001;
- U.S. Patent Application Serial No. _____, entitled “Derivation and Quantization of Robust Non-Local Characteristics for Blind Watermarking” filed on April 24, 2001;
- U.S. Patent Application Serial No. _____, entitled “Recognizer of Audio-Content in Digital Signals” filed on April 24, 2001
- U.S. Patent Application Serial No. 09/259,669, entitled “A System and Method for Producing Modulated Complex Lapped Transforms” filed on February 26, 1999; and
- U.S. Patent Application Serial No. 09/421,986, entitled “System and Method for Hashing Digital Images” filed on October 19, 1999.

Introduction

Described herein are one or more exemplary implementations of a recognizer of text-based work. One of such exemplary implementations may be referred to as an exemplary “text recognizer.”

One implementation of an exemplary text recognizer described herein, automatically and accurately detects plagiarism in text-based work based upon the text content of such work. Another implementation of an exemplary text recognizer, described herein, automatically and accurately categorizes text-based work based upon the text content of such work.

These exemplary implementations may be implemented (whole or in part) by a text recognition system 100 and/or by a computing environment like that shown in Fig. 7.

Brief Overview

An exemplary text recognizer determines one or more recognition representations (e.g., hash values) for the contents of a text document. These representations may include an intermediate and a final hash value. In general, it does this by:

- filtering out non-essential words, punctuation, and such;
- putting it into a standard format;
- extracting sub-text via a self-synchronized approach, such as:
 - fixed-length subtext extraction; or
 - variable-length subtext extraction;

- arranging the extracted sub-text into a standard format (i.e., an image)
- hashing the sub-text or image

One of the results of hashing the image (or the sub-text itself) is a final hash value, which uniquely identifies the body of the text document. This hash value does not rely on any particular order of text, any punctuation, any non-essential words. In addition, the selected words (i.e., sub-text) that are the basis for the hash value calculation are chosen via a pseudorandom and/or cryptographic fashion; thus, making it exceptionally difficult for a plagiarizer to predict which words are essential.

Furthermore, another result of hashing the image is an intermediate hash value, which generally identifies the body of the text document. Bodies of text having semantically similar content have similar intermediate hash values. Thus, their hash values cluster together.

Perceptually Same and Perceptually Distinct

The exemplary text recognizer treats two “perceptually same” bodies of text as the same if a human observer reasonably views them as the same. This may also be called “perceptually identical,” “imperceptibly indistinguishable,” or other similar equivalent phrases. For example, “perceptually same” bodies of text are documents that look as if they are they are substantially same to a human.

In contrast, a “perceptually distinct” digital goods is generally the converse of “perceptually same” digital goods. This may also be called “perceptually different” or “perceptually distinguishable”.

Perceptually Similar

The exemplary text recognizer treats two “perceptually similar” bodies of text as documents that should be categorized as similar. Herein, a pair of bodies of text are “perceptually similar” when their “similarity-recognition” values are close in value (i.e., proximal). In other words, the phrase “perceptually similar” stresses the fact that two documents should produce the proximally near hash values.

Exemplary Text Recognition System

Figs. 1-3 show a text identification sub-system 100, a similarity detection (alternatively, a plagiarism detection) sub-system 200, and a text categorization sub-system 300. The sub-systems alone or in combination with each other form an embodiment of an exemplary text recognizer. Each embodiment may be implemented in software, hardware, or a combination thereof.

Text Identification Sub-System

Fig. 1 shows the text identification sub-system 100, which is an embodiment (or a portion thereof) of the exemplary text recognizer. This embodiment includes a text-based work retriever 110 for retrieving text-based works from a database of text-based works, such as database 112, or some other source of text-based works.

A text-based work may also be called a “body of text”, a “text document”, “text file”, or other similar labels that refers to an object having content consisting

1 primarily of text. In particular, the text is electronically encoded in a manner
2 readable by a computer.

3 Once a subject work is chosen, a text filter 120 filters out the non-essential
4 (i.e., superfluous) elements, which includes non-essential and/or non-distinctive
5 words, punctuations, symbols, and the like. For a given application, the definition
6 of non-essential elements may be customized. Examples of non-essential elements
7 include common words and virtually everything that is not a word. Examples of
8 common words include:

- 9 • articles (e.g., "a", "an", "the");
- 10 • pronouns (e.g., "I", "he", "she");
- 11 • prepositions (e.g., "of", "on", "for");
- 12 • conjunctions (e.g., "and", "or", "but");
- 13 • common verbs (e.g., "is", "are", "were"); and
- 14 • other such common words.

15 After filtering the original text, a text formatter 130 formats the remaining
16 text into a canonical (i.e., standard) format. For example, the canonical format
17 may ignore all characters that are not letters or spaces. In the canonical format all
18 of the letters may be converted to uppercase.

19 An example of some original text may be as follows:

20
21 A text-based work may also be called a
22 "body of text", a "text document", "text
23 file", or other similar labels that refers to
24 an object having content consisting primarily
25 of text. In particular, the text is

1 electronically encoded in a manner readable by
2 a computer.

3
4 One possible example of canonical, filtered text that could result from the
5 above example (filtered out common words, converted to all caps, no characters
6 but letters and spaces) is:

7
8 TEXT BASED WORK MAY ALSO CALLED BODY TEXT
9 TEXT DOCUMENT TEXT FILE OTHER SIMILAR LABELS
10 REFERS OBJECT HAVING CONTENT CONSISTING
11 PRIMARILY TEXT PARTICULAR TEXT ELECTRONICALLY
12 ENCODED MANNER READABLE COMPUTER
13

14 A sub-text extractor 140 extracts the selected sub-text from the canonically
15 formatted text. The extractor 140 uses a cryptographic key to guide its selection of
16 sub-text from the canonically formatted text. This uses conventional cryptographic
17 techniques, such as a pseudo-random number generator.

18 This extraction may be accomplished using a fix-length approach or a
19 variable-length approach. Both of these approaches are self-synchronized. This
20 extraction technique (and these two approaches) is described in detail later in the
21 Text Sifting section.

22 A sub-text imager 150 may arrange the extracted sub-text into a standard
23 format called an "image." For example, the "image" may be literally an image of
24 the extracted sub-text within a standard format, standard font, standard
25 background, etc. The text in the image may be black against a white background.

1 Alternatively, the text and background may have different colors and different
2 contrasts.

3 The standard format may also be called “formatted sifted text.” A hasher
4 160 finds hash values of such an image. The resulting hash values of the hasher
5 160 include an intermediate hash value (i.e., categorization hash value) and a final
6 hash value. These hash values are recognition representations of the text of the
7 original body of text. That is because these hash values may be used to recognize
8 (and even identify) the text within a body of text.

9 10 Hashing

11 Hashing techniques have been used to protect the rights of content owners.
12 In general, a hashing technique maps a large block of raw data into relatively small
13 and structured set of identifiers. These identifiers are also referred to as “hash
14 values” or simply “hash.” By introducing a specific structure and order into raw
15 data, the hashing function drastically reduces the size of the raw data into short
16 identifiers. It simplifies many data management issues and reduces the
17 computational resources needed for accessing large databases.

18 Mathematically, a hashing technique include a hashing function $H_K(\cdot)$. That
19 function takes a signal x as input and computes a short vector $h = H_K(x)$. That
20 vector is an apparently random value, which is indexed by a secret key K , in some
21 large set. That vector h is a hash value.

22 The hash values produced by such techniques are viewed as useful because
23 they typically have following desirable characteristics:
24
25

- 1 • Apparently Uniformly Distributed—For any given input, the output
2 hash value are uniformly distributed among the possible L -bit
3 outputs.
- 4 • Approximate Pairwise Independent—For two distinct inputs, the
5 corresponding outputs are statistically almost independent of each
6 other.

7 Additional details regarding calculating hash values are discussed in the
8 pending U.S. Patent Applications that are incorporated by reference.

9 [kcc1]The resulting hash values may be displayed and stored. These values
10 are stored in the database 112 (or some other database) and associated with the
11 original suspect work from which the values were calculated.

12 Alternatively, the text identification sub-system 100 may skip hasher 160
13 and store the sub-text image (i.e., formatted sifted text) in database 112 in
14 association with the original work. This sub-text image may be used later to
15 compare with the sub-text image of another work to detect similarity (e.g.,
16 plagiarism) (or to categorize the work).

17 Similarity Detection Sub-System

18
19 Fig. 2 shows the similarity detection sub-system 200, which is an
20 embodiment (or a portion thereof) of the exemplary text recognizer. This
21 embodiment includes a hash-value retriever 210 for retrieving a hash-value of a
22 selected text-based work. More particularly, it retrieves the final hash value of
23 such work from a database of text-based works, such as the database 112 (shown
24 in Figs. 1-3), or some other source of text-based works.

1 Fig. 2 also shows the text identification sub-system 100 (shown in Fig.1
2 and described above). It calculates a final hash value of a suspect text-base work
3 222. This work is one that is suspected of plagiarizing a work found in a database
4 of text-based works—such as the database 112 (shown in Figs. 1-3) or some other
5 source of text-based works.

6 The text identification sub-system 100 provides the final hash value of the
7 suspect work 222 to a hash value comparator 230. Likewise, the retriever 210
8 provides the final hash value of the selected work to the comparator 230. Of
9 course, this can be reversed so that the sub-system 100 provides the hash value of
10 the *selected* work while the retriever 210 provides the hash value of the *suspected*
11 work.

12 As its name suggests, the comparator 230 compares the hash values of the
13 two works to determine if they substantially match. Substantially matching means
14 that the two hash values are close enough in value to reasonably conclude that the
15 two works have the same hash values within a margin of error. This margin of
16 error may be subjectively determined by a user or designer for a specific
17 application.

18 The results of the comparison are displayed on a display 250 and stored in
19 a data store 250. The results indicate whether the content of the suspect work 222
20 is plagiarized from the selected work.

21 Alternatively, the similarity detection sub-system 200 may use the sub-text
22 itself or the sub-text image (i.e., formatted sifted text) of a work rather than its
23 hash value. With this alternative embodiment, the sub-text images of the work
24 may be compared or the hash values of the works may be calculated and those
25 values compared.

Text Categorization Sub-System

Fig. 3 shows the text categorization sub-system 300, which is an embodiment (or a portion thereof) of the exemplary text recognizer. This embodiment includes a hash-value retriever 210 for retrieving a hash-value of a selected text-based work. More particularly, it retrieves the intermediate (i.e., categorization) hash value of such work from a database of text-based works, such as the database 112 (shown in Figs. 1-3), or some other source of text-based works. Of course, the retriever 210 could retrieve an intermediate hash value of a work processed by the text identifier sub-system 100.

A work categorizer 330 uses the categorization hash value of the selected work to group such work with others of similar (i.e., proximal) categorization hash values. In other words, based upon the categorization hash value of a given work, the work categorizer 330 groups the given work with other works having similar categorization hash values. Thus, the hash values of all works in each grouping are clustered together (i.e., proximal each other). Although these groupings are objectively determined, the subjective nature of the content of works within a grouping will be similar to that of the content of others within the grouping.

The boundaries between groupings are determined manually or automatically. Manually, a person selects the boundary between groupings using the natural clustering seen after many works have been categorized. Automatically, a system mathematically selects the boundary between groupings to be some point between (perhaps halfway) the centers of the groupings. Of course, other such techniques may be used to determine boundaries. These techniques may be fully automatic, fully manual, or some combination.

1 The work categorizer 330 displays the results on a display 350 and stores it
2 in a data store 340 and/or database 112.

3 Alternatively, the text categorization sub-system 300 may use the sub-text
4 itself or the sub-text image (i.e., formatted sifted text) of a work rather than its
5 hash value. With this alternative embodiment, the categorization hash value of a
6 work is calculated and the result is used to categorize the work.

7 **Methodological Implementation of the Exemplary Text Recognizer**

9 Fig. 4 shows an illustrative methodological implementation of the
10 exemplary text recognizer performed by the text identification sub-system 100, the
11 similarity detection sub-system 200, or the text categorization sub-system 300 (or
12 some combination of such sub-systems). This methodological implementation
13 may be performed in software, hardware, or a combination thereof.

14 **Text Identification Methodological Implementation**

16 Fig. 4 illustrates a text identification methodological implementation of the
17 exemplary text recognizer. At 410 of Fig. 4, the exemplary text recognizer
18 retrieves a subject text-based work from a database of text-based works or some
19 other source of such works. Once a subject work is chosen, the exemplary text
20 recognizer, at 412, filters out the non-essential (i.e., superfluous) elements, which
21 includes non-essential and non-distinctive words, punctuations, symbols, and the
22 like. Examples are provided above in the discussion of text filter 120 of Fig. 1.

23 At 414 of Fig. 4, the exemplary text recognizer formats the remaining text
24 into a canonical (i.e., standard) format. At 416, it extracts selected sub-text from
25 the canonically formatted text. The extraction uses a cryptographic key to guide its

1 selection of sub-text from the canonically formatted text. This extraction may be
2 accomplished using a fix-length approach or a variable-length approach. This
3 extraction technique (and these two approaches) is described in detail later in the
4 Text Sifting section.

5 At 418 of Fig. 4, the exemplary text recognizer arranges the extracted sub-
6 text into a standard format called an "image." At 420, it calculates a hash value of
7 such image. The hash values include an intermediate hash value (i.e.,
8 categorization hash value) and a final hash value. These hash values are
9 recognition representations of the text of the original body of text. That is because
10 these hash values may be used to recognize (and even identify) the text within a
11 body of text.

12 At 422 of Fig. 4, the resulting hash values are displayed and stored. These
13 values are stored in a database in association with the original subject work from
14 which the values were calculated. These hash values may be used later to compare
15 with the hash values of another work to detect similarity (or to categorize the
16 work).

17 Alternatively, the text identification methodological implementation may
18 skip the hash value calculation block 420 and store, at block 422, the sub-text
19 itself or the sub-text image (i.e., formatted sifted text) in association with the
20 original work. This sub-text image may be used later to compare with the sub-text
21 image of another work to detect similarity (e.g., plagiarism) (or to categorize the
22 work).

Similarity Detection Methodological Implementation

Fig. 5 illustrates a similarity detection methodological implementation of the exemplary text recognizer. At 456 of Fig. 5, the exemplary text recognizer retrieves a hash value of a selected text-based work. More particularly, it retrieves the final hash value of such work from a database of text-based works or some other source of such works.

Fig. 5 also shows the text identification method of Fig. 4 at block 452. The method 452 calculates a final hash value of a text-based work 450 that is suspected of being copied from the selected work retrieved by block 456. At 454, the exemplary text recognizer retrieves the calculated final hash value of the suspect text-base work 450 from the text identification method of block 452. Of course, this can be reversed so that the method 452 provides the hash value of the *selected* work while block 452 provides the hash value of the *suspected* work.

At 458, the exemplary text recognizer compares the hash values of the two works (suspect work 450 and selected work of 456) to determine if they substantially match. Substantially matching means that the two hash values are close enough in value to reasonably conclude that the two works have the same hash values within a margin of error.

If the result of such comparison is no substantial match, then the exemplary text recognizer indicates, at 460, that the suspect work 450 is not a substantial copy of the selected work of 456. In other words, no plagiarism is detected if the final hash values of compared works do not substantially match. At 464, this process ends.

1 However, if the result of such comparison is a substantial match, then the
2 exemplary text recognizer indicates, at 462, that the suspect work 450 is a
3 substantial copy of the selected work of 456. In other words, plagiarism is
4 detected if the final hash values of compared works substantially match. At 464,
5 this process ends.

6 Alternatively, the similarity detection methodological implementation may
7 use the sub-text itself or the sub-text image (i.e., formatted sifted text) of a work
8 rather than its hash value. With this alternative embodiment, the sub-text images
9 of work may be compared or the hash values of the works may be calculated and
10 those values compared.

11 Text Categorization Methodological Implementation

12
13 Fig. 6 illustrates a text categorization methodological implementation of the
14 exemplary text recognizer. At 516 of Fig. 6, the exemplary text recognizer
15 retrieves a hash value of a selected text-based work. More particularly, it retrieves
16 the intermediate (i.e., categorization) hash value of such work from a database of
17 text-based works or some other source of such works.

18 In dashed box 505, Fig. 6 also shows an alternative way of getting an
19 intermediate hash value of the selected work. This is by processing the work using
20 the text identification method of Fig. 4 at block 512. The method 512 calculates
21 an intermediate (i.e., categorization) hash value of the selected work. At 514, the
22 exemplary text recognizer retrieves the calculated intermediate hash value of the
23 selected text-base work 510 from the text identification method of block 512.

24 At 520, the exemplary text recognizer uses the intermediate hash value of
25 the selected work to group such work with others of similar (i.e., proximal)

intermediate hash values. In other words, based upon the intermediate hash value of a given work, the exemplary text recognizer groups the given work with other works having similar intermediate hash values. Thus, the hash values of all works in a given grouping are clustered together (i.e., proximal each other). Although these groupings are objectively determined, the subjective nature of the content of works within a grouping will be similar to that of the content of others within the grouping.

See the above description of the text categorization sub-system (of Fig. 3) to see how the boundaries between groupings may be determined.

At 522, the exemplary text recognizer stores the categorization results in a database. At 524, the process ends.

Alternatively, the text categorization methodological implementation may use the sub-text itself or the sub-text image (i.e., formatted sifted text) of a work rather than its hash value. With this alternative embodiment, the categorization hash value of a work is calculated and the result is used to categorize the work.

Text Sifting

A technique called “cryptographic text sifting” (or simply “text sifting” herein) takes a large body of text and characterizes it in a simpler and identifying form. The characterization may be by a text subset selected from such body of text. From this subset, one or more hash values may be calculated.

Text-sifting technique is performed, wholly or in part, by the text identification sub-system 100 shown in Fig. 1. This includes: the text filter 120, the text formatter 130, the sub-text extractor 140, the sub-text imager 150, and the

1 hasher 160. Furthermore, the text-sifting technique is performed, wholly or in
2 part, by the text identification methodological implementation illustrated in Fig. 4.

3 The cryptographic text-sifting technique takes as input a text document
4 (i.e., a text-based work) and a secret key. It selects and outputs a small number of
5 words from the document. Without access to the secret key, an adversary (e.g., a
6 plagiarist) cannot predict which words will be selected. The adversary can only
7 affect the output by making changes to document.

8 This helps immensely in the effort to detect copyright infringement.
9 Instead of combing through an entire database to see if a file has been plagiarized,
10 the database might contain sifting versions of all documents alongside the
11 complete versions. Searching for matches in the sifted versions is much faster than
12 searching through the complete files. Alternatively, hash values of such sifted
13 version may be stored instead of and/or in addition to the sifted versions
14 themselves. Consequently, the hash values may be compared instead of the sifted
15 versions. Furthermore, if sifted versions are stored, hash values can be calculated
16 and then compared.

17 If a substantial match is found, the complete versions can then be
18 compared. If an adversary wanted to change a plagiarized document so much that
19 the sifting versions had nothing common, the adversary would need to change so
20 much of the document that a person is unlikely to recognize the original.

21 These text-sifted versions of the original text are recognition
22 representations of the text of the original body of text. That is because they may
23 be used to recognize (and even identify) the text within a body of text.

24 There are two main approaches to text sifting. Both are self-synchronizing.
25 The first approach is *constant length text sifting*. It may also be called the *fixed-*

length approach. With this approach, the number of words outputted by the sifting technique is constant. Second approach is *proportional length text sifting*. It may also be called the *variable-length approach*. With this approach, the sifting technique outputs a number of words proportional to the length of the document.

General Introduction to Text Sifting

For the exemplary text recognizer, the text-sifting technique takes as input a text document (i.e., a text-based work) and a secret key. It outputs a subset of words from that document. This subset may be further formatted into a standard “image,” which may be called “formatted sifted text.”

The text-sifting technique uses hash functions to decide which words will be selected. The specific words that are chosen are always the same when the same document and key are used. An adversary, who has access to the document, but not the secret key, cannot predict which words will be selected.

Groups of m Words

In the exemplary text recognizer, text-sifting technique works on groups of m words (m -tuples of words) instead of only single words. Those groups of words need not be contiguous. For example, consider the following sentence:

She sells seashells by the seashore

The sentence consists of the following 2-tuples:

She sells

Sells seashells

seashells by

by the
the seashore

It also consists of the following 3-tuples:

she sells seashells
sells seashells by
seashells by the
by the seashore

If a document has N words, there are $(N - m)$ m -tuples in the document. Since m is small compared to N , n is approximately equal to N . Excluding the words at the very beginning or the end of a document, each word is part of m of the m -tuples. This difference is minor, again because m is small compared to N .

There several reasons for operating on m -tuples instead of single words. First, since m can be set equal to one, this does not preclude single words, and makes the text-sifting technique more general. Second, m -tuples of words are more easily distinguishable than are single words, which is useful in applications of text sifting (such as similarity detection). Finally, using m -tuples in effect links a word to his neighbors, which helps the text-sifting technique take a word's place in a document into account.

In addition, different sizes of m -tuples may be used in the same document. For example, a document may be evaluated several times, first with $m=3$ then $m=5$ then $m=15$, and taking the results of all three times into account.

An alternative approach to m -tuples is "windowing." With this approach, a window onto the text is used and a specific number of words are selected from that

1 window. For example, a window may be twenty words long and only seven words
2 are selected from that window.

3 4 Removal of Superfluous Details

5 In addition to working on groups of m words, the exemplary text recognizer
6 ignores superfluous details that will not distinguish documents from each other.
7 Such superfluous details may be called non-essential (i.e., superfluous) elements.

8 This is accomplished in two separate ways. First, before any text in a
9 document is sifting, all white space is converted to single spaces, all characters
10 that are not letters or spaces are purged from the document, and all letters are
11 converted to uppercase. Second, all words found in a list of common words are
12 removed from the document. This removal of superfluous details helps insure that
13 the words that are selected deal with the substantive intellectual content of the
14 document, rather than the formatting or unimportant words that can easily
15 changed.

16 In addition, the format (i.e., type) of document may be identified. Examples
17 of formats include HTML, postscript, ASCII, etc. Doing this, the document may
18 be analyzed and customize according to its identified format so that non-essential
19 elements may be removed intelligently.

20 For a postscript document, this would consist of removing the postscript
21 commands from the document, and likewise it would remove the HTML
22 commands from an HTML document. For other text-based documents, it could
23 remove all of the formatting characters, etc.

1 If removal of non-essential elements is not conditioned on the format of the
2 documents, then a match might be made simply based upon formatting
3 similarities.

4 Further Defining Text-Sifting Technique

6 A cryptographic text-sifting technique f takes as input a secret key K and a
7 document D consisting of n ($= N - m$) m -tuples, $w_1, w_2, w_3, \dots, w_n$. The output $f(D,$
8 $K)$ is S , a subset of D that consists of k of the w_i . For an adversary without K , the k
9 m -tuples in S appear to be chosen uniformly at random from all m -tuples in D .

10 Note that a specific word in D might be selected more than once. For
11 example, consider the previous example. If 2-tuples are being used, the text-
12 sifting technique could use "she sells" and "sells seashells," causing the word
13 "sells" to appear twice.

14 In addition, a specific m -tuple may appear more than once.

15 Proportional-Length Sifting Approach

17 The proportional-length sifting approach may also be called the variable-
18 length approach. This approach generates output of a size proportional to the size
19 of input (the original document). The size of the output is on average proportional
20 to the size of the input. That is, the expected value of the size is proportional to
21 the size of the input. The larger the document is, the less variation there is in this
22 proportionality.

23 A proportional-length cryptographic text-sifting technique is a
24 cryptographic text-sifting technique where k is proportional to n for all documents
25 D . $k = cn$ for some constant c , such that $1/c$ is an integer.

To achieve a proportional length text-sifting technique, the exemplary text recognizer hashes each m -tuple and looks at the value modulo $(1/c)$. For example, if a proportional length text-sifting technique should output one m -tuple for each 20 m -tuples in D , then $c = .05$ and an m -tuple will be selected only if its hash value is congruent to 0 mod 20.

If $1/c$ is not an integer, then accept an m -tuple if its hash is less than $x \bmod y$, where $c = x/y$.

A hash function is used to determine which words will be selected. This unrelated to the final and intermediate hash values calculated by the hasher 160 and by block 430. To distinguish them, the hash values used to determine which words would be selected is called the "selected hash values."

To compute the selected hash value for each word, the hash function is performed. It takes as input the secret key k and an m -tuple and outputs a b -bit number. It treats the m -tuple as a number by considering each character to be a digit in a base corresponding to a range of characters. For example, the range may be 0 to 26 for the letters and the space character) and therefore a digit in a base-27 number. The m -tuple corresponds to this number, called x . The function also computes k random coefficients, chosen uniformly and the range of $[0, 2^b]$, called $c_1, c_2, c_3, \dots, c_k$. Finally, the function chooses a random b -bit prime number p . The select hash value, h , is:

$$h = c_0x^k + c_1x^{k-1} + c_2x^{k-2} + \dots + c_{k-2}x^2 + c_{k-1}x^1 + c_k \quad \bmod p \bmod (1/c)$$

To speed up the process, Horner's method is used to evaluate h . Horner's method iteratively re-evaluates h in place as follows:

$$\begin{array}{llll}
 \text{Step 1} & h = c_0 & & \text{mod } p \\
 \text{Step 2} & h = hx + c_1 = c_0x + c_1 & & \text{mod } p \\
 \text{Step 3} & h = hx + c_2 = c_0x^2 + c_1x + c_2 & & \text{mod } p \\
 \text{Step 4} & h = hx + c_3 = c_0x^3 + c_1x^2 + c_2x + c_3 & & \text{mod } p \\
 & : & & : \\
 \text{Step } k-1 & h = hx + c_{k-1} = c_0x^{k-1} + c_1x^{k-2} + c_2x^{k-3} + \dots + c_{k-3}x^2 + c_{k-2}x + c_{k-1} & & \text{mod } p \\
 \text{Step } k & h = hx + c_k = c_0x^k + c_1x^{k-1} + c_2x^{k-2} + \dots + c_{k-2}x^2 + c_{k-1}x + c_k & & \text{mod } p
 \end{array}$$

A final value of h is taken mod $(1/c)$ and if the result is zero, is output.

Alternatively, this may be described as follows: c is a rational number of the form $c=a/b$, where both a and b are integers. An m -tuple is output if its hash value is less than $a \bmod b$. Another way of saying it: an m -tuple is output if it is $0, 1, 2, \dots, a-2 \bmod b$.

Constant-length Sifting Approach The constant-length sifting approach may also be called the fixed-length approach. It may be desirable to have a text-sifting technique that always outputs a constant number of m -tuples. The constant-length sifting approach achieves this. With the constant-length cryptographic text-sifting technique, k is fixed for all documents D .

To do this, the approach uses the secret key to hash every m -tuple in the document. It selects the m -tuple with the smallest hash value. Alternatively, it could select the largest hash value or some other relative measure.

1 To expand this approach to output k m -tuples, the exemplary text
2 recognizer simply repeats the function k times with k different hash functions.

3 Like the previously discussed proportional-length sifting approach, a hash
4 function is used to determine which words will be selected in the constant-length
5 sifting approach. This unrelated to the final and intermediate hash values
6 calculated by the hasher 160 and by block 430. To distinguish them, the hash
7 values used to determine which words would be selected is called the “selected
8 hash values.”

9 It may be very slow to hash every m -tuple with k different hash functions;
10 so two strategies are employed for achieving a reasonable level of speed. The first
11 is to use a very fast hash function. The second is to correlate the hash functions so
12 that computing k select hash values at once is faster the computing k separate
13 select hash values.

14 To generate all these select hash values, $h_1, h_2, h_3, \dots, h_k$, a single task
15 computes several hash functions. The single task takes as input the secret key k
16 and an m -tuple and outputs a b -bit number. It treats the m -tuple as a number, x ,
17 just as the proportional-length sifting approach. The task also computes $2k-1$
18 random coefficients, chosen uniformly the range of $[0, 2^b]$, called $c_1, c_2, c_3, \dots, c_{2k-1}$.
19 Finally, the task chooses a random b -bit prime number p .

To compute the select hash values, the following is done:

$$h_1 = c_0x^k + c_1x^{k-1} + c_2x^{k-2} + \dots + c_{k-2}x^2 + c_{k-1}x^1 + c_k \quad \text{mod } p$$

$$h_2 = c_1x^k + c_2x^{k-1} + c_3x^{k-2} + \dots + c_{k-1}x^2 + c_kx^1 + c_{k+1} \quad \text{mod } p$$

$$h_3 = c_2x^k + c_3x^{k-1} + c_4x^{k-2} + \dots + c_kx^2 + c_{k+1}x^1 + c_{k+2} \quad \text{mod } p$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \text{mod } p$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad \text{mod } p$$

$$h_{k-1} = c_{k-1}x^k + c_kx^{k-1} + c_{k+1}x^{k-2} + \dots + c_{2k-4}x^2 + c_{2k-3}x^1 + c_{2k-2} \quad \text{mod } p$$

$$h_k = c_kx^k + c_{k+1}x^{k-1} + c_{k+2}x^{k-2} + \dots + c_{2k-3}x^2 + c_{2k-2}x^1 + c_{2k-1} \quad \text{mod } p$$

The first step of this process is to compute $x^2 \text{ mod } p$, $x^3 \text{ mod } p$, $x^4 \text{ mod } p$, ..., $x^k \text{ mod } p$. Next, the x^i are multiplied by c_i and summed modulo p to form h_1 as above. Thus, computing h_1 requires $(2k-1)$ modular multiplications and k modular additions.

Computing each successive h_i is done according to following formula:

$$h_{i+1} = h_i x + c_{i+k-1} - (c_{i-1}x^k) \quad \text{mod } p$$

Since x^k was already calculated, each successive h_i takes 2 modular multiplications, one modular addition, and one modular subtraction.

$$\text{Total cost} = (\text{cost of computing } h_1) + (k-1)(\text{cost of computing successive } h_i)$$

Self-Synchronization

Both of these approaches are self-synchronizing. This means that it outputs words (or *m*-tuples) based on their own inherent features, not based on other features of the document. For example, adding a word at the beginning of the document will not completely change all of the output. As are result, the body of text is less sensitive to de-synchronizing attacks (e.g., scrambling, rearranging, deletions, insertions, etc.)

Exemplary Computing System and Environment

Fig. 7 illustrates an example of a suitable computing environment 900 within which an exemplary text recognizer, as described herein, may be implemented (either fully or partially). The computing environment 900 may be utilized in the computer and network architectures described herein.

The exemplary computing environment 900 is only one example of a computing environment and is not intended to suggest any limitation as to the scope of use or functionality of the computer and network architectures. Neither should the computing environment 900 be interpreted as having any dependency or requirement relating to any one or combination of components illustrated in the exemplary computing environment 900.

The exemplary text recognizer may be implemented with numerous other general purpose or special purpose computing system environments or configurations. Examples of well known computing systems, environments, and/or configurations that may be suitable for use include, but are not limited to, personal computers, server computers, thin clients, thick clients, hand-held or

1 laptop devices, multiprocessor systems, microprocessor-based systems, set top
2 boxes, programmable consumer electronics, network PCs, minicomputers,
3 mainframe computers, distributed computing environments that include any of the
4 above systems or devices, and the like.

5 The exemplary text recognizer may be described in the general context of
6 computer-executable instructions, such as program modules, being executed by a
7 computer. Generally, program modules include routines, programs, objects,
8 components, data structures, etc. that perform particular tasks or implement
9 particular abstract data types. The exemplary text recognizer may also be
10 practiced in distributed computing environments where tasks are performed by
11 remote processing devices that are linked through a communications network. In
12 a distributed computing environment, program modules may be located in both
13 local and remote computer storage media including memory storage devices.

14 The computing environment 900 includes a general-purpose computing
15 device in the form of a computer 902. The components of computer 902 can
16 include, by are not limited to, one or more processors or processing units 904, a
17 system memory 906, and a system bus 908 that couples various system
18 components including the processor 904 to the system memory 906.

19 The system bus 908 represents one or more of any of several types of bus
20 structures, including a memory bus or memory controller, a peripheral bus, an
21 accelerated graphics port, and a processor or local bus using any of a variety of
22 bus architectures. By way of example, such architectures can include an Industry
23 Standard Architecture (ISA) bus, a Micro Channel Architecture (MCA) bus, an
24 Enhanced ISA (EISA) bus, a Video Electronics Standards Association (VESA)
25

1 local bus, and a Peripheral Component Interconnects (PCI) bus also known as a
2 Mezzanine bus.

3 Computer 902 typically includes a variety of computer readable media.
4 Such media can be any available media that is accessible by computer 902 and
5 includes both volatile and non-volatile media, removable and non-removable
6 media.

7 The system memory 906 includes computer readable media in the form of
8 volatile memory, such as random access memory (RAM) 910, and/or non-volatile
9 memory, such as read only memory (ROM) 912. A basic input/output system
10 (BIOS) 914, containing the basic routines that help to transfer information
11 between elements within computer 902, such as during start-up, is stored in ROM
12 912. RAM 910 typically contains data and/or program modules that are
13 immediately accessible to and/or presently operated on by the processing unit 904.

14 Computer 902 may also include other removable/non-removable,
15 volatile/non-volatile computer storage media. By way of example, Fig. 7
16 illustrates a hard disk drive 916 for reading from and writing to a non-removable,
17 non-volatile magnetic media (not shown), a magnetic disk drive 918 for reading
18 from and writing to a removable, non-volatile magnetic disk 920 (e.g., a "floppy
19 disk"), and an optical disk drive 922 for reading from and/or writing to a
20 removable, non-volatile optical disk 924 such as a CD-ROM, DVD-ROM, or other
21 optical media. The hard disk drive 916, magnetic disk drive 918, and optical disk
22 drive 922 are each connected to the system bus 908 by one or more data media
23 interfaces 926. Alternatively, the hard disk drive 916, magnetic disk drive 918,
24 and optical disk drive 922 can be connected to the system bus 908 by one or more
25 interfaces (not shown).

1 The disk drives and their associated computer-readable media provide non-
2 volatile storage of computer readable instructions, data structures, program
3 modules, and other data for computer 902. Although the example illustrates a hard
4 disk 916, a removable magnetic disk 920, and a removable optical disk 924, it is to
5 be appreciated that other types of computer readable media which can store data
6 that is accessible by a computer, such as magnetic cassettes or other magnetic
7 storage devices, flash memory cards, CD-ROM, digital versatile disks (DVD) or
8 other optical storage, random access memories (RAM), read only memories
9 (ROM), electrically erasable programmable read-only memory (EEPROM), and
10 the like, can also be utilized to implement the exemplary computing system and
11 environment.

12 Any number of program modules can be stored on the hard disk 916,
13 magnetic disk 920, optical disk 924, ROM 912, and/or RAM 910, including by
14 way of example, an operating system 926, one or more application programs 928,
15 other program modules 930, and program data 932. Each of such operating
16 system 926, one or more application programs 928, other program modules 930,
17 and program data 932 (or some combination thereof) may include an embodiment
18 of text-based work retriever, text filter, text formatter, sub-text extractor, sub-text
19 imager, hasher, text sifter, text identification sub-system, similarity detection sub-
20 system, categorization sub-system, recognition representation determiner, hash
21 value comparator, hash value retriever, text works database, and work categorizer.

22 A user can enter commands and information into computer 902 via input
23 devices such as a keyboard 934 and a pointing device 936 (e.g., a "mouse").
24 Other input devices 938 (not shown specifically) may include a microphone,
25 joystick, game pad, satellite dish, serial port, scanner, and/or the like. These and

1 other input devices are connected to the processing unit 904 via input/output
2 interfaces 940 that are coupled to the system bus 908, but may be connected by
3 other interface and bus structures, such as a parallel port, game port, or a universal
4 serial bus (USB).

5 A monitor 942 or other type of display device can also be connected to the
6 system bus 908 via an interface, such as a video adapter 944. In addition to the
7 monitor 942, other output peripheral devices can include components such as
8 speakers (not shown) and a printer 946 which can be connected to computer 902
9 via the input/output interfaces 940.

10 Computer 902 can operate in a networked environment using logical
11 connections to one or more remote computers, such as a remote computing device
12 948. By way of example, the remote computing device 948 can be a personal
13 computer, portable computer, a server, a router, a network computer, a peer device
14 or other common network node, and the like. The remote computing device 948 is
15 illustrated as a portable computer that can include many or all of the elements and
16 features described herein relative to computer 902.

17 Logical connections between computer 902 and the remote computer 948
18 are depicted as a local area network (LAN) 950 and a general wide area network
19 (WAN) 952. Such networking environments are commonplace in offices,
20 enterprise-wide computer networks, intranets, and the Internet.

21 When implemented in a LAN networking environment, the computer 902 is
22 connected to a local network 950 via a network interface or adapter 954. When
23 implemented in a WAN networking environment, the computer 902 typically
24 includes a modem 956 or other means for establishing communications over the
25 wide network 952. The modem 956, which can be internal or external to computer

1 902, can be connected to the system bus 908 via the input/output interfaces 940 or
2 other appropriate mechanisms. It is to be appreciated that the illustrated network
3 connections are exemplary and that other means of establishing communication
4 link(s) between the computers 902 and 948 can be employed.

5 In a networked environment, such as that illustrated with computing
6 environment 900, program modules depicted relative to the computer 902, or
7 portions thereof, may be stored in a remote memory storage device. By way of
8 example, remote application programs 958 reside on a memory device of remote
9 computer 948. For purposes of illustration, application programs and other
10 executable program components such as the operating system are illustrated herein
11 as discrete blocks, although it is recognized that such programs and components
12 reside at various times in different storage components of the computing device
13 902, and are executed by the data processor(s) of the computer.

14 **Computer-Executable Instructions**

15
16 An implementation of an exemplary text recognizer may be described in the
17 general context of computer-executable instructions, such as program modules,
18 executed by one or more computers or other devices. Generally, program modules
19 include routines, programs, objects, components, data structures, etc. that perform
20 particular tasks or implement particular abstract data types. Typically, the
21 functionality of the program modules may be combined or distributed as desired in
22 various embodiments.

Exemplary Operating Environment

Fig. 7 illustrates an example of a suitable operating environment 900 in which an exemplary text recognizer may be implemented. Specifically, the exemplary text recognizer(s) described herein may be implemented (wholly or in part) by any program modules 928-930 and/or operating system 926 in Fig. exemplary text recognizer or a portion thereof.

The operating environment is only an example of a suitable operating environment and is not intended to suggest any limitation as to the scope or use of functionality of the exemplary text recognizer(s) described herein. Other well known computing systems, environments, and/or configurations that are suitable for use include, but are not limited to, personal computers (PCs), server computers, hand-held or laptop devices, multiprocessor systems, microprocessor-based systems, programmable consumer electronics, wireless phones and equipments, general- and special-purpose appliances, application-specific integrated circuits (ASICs), network PCs, minicomputers, mainframe computers, distributed computing environments that include any of the above systems or devices, and the like.

Computer Readable Media

An implementation of an exemplary text recognizer may be stored on or transmitted across some form of computer readable media. Computer readable media can be any available media that can be accessed by a computer. By way of example, and not limitation, computer readable media may comprise "computer storage media" and "communications media."

1 “Computer storage media” include volatile and non-volatile, removable and
2 non-removable media implemented in any method or technology for storage of
3 information such as computer readable instructions, data structures, program
4 modules, or other data. Computer storage media includes, but is not limited to,
5 RAM, ROM, EEPROM, flash memory or other memory technology, CD-ROM,
6 digital versatile disks (DVD) or other optical storage, magnetic cassettes, magnetic
7 tape, magnetic disk storage or other magnetic storage devices, or any other
8 medium which can be used to store the desired information and which can be
9 accessed by a computer.

10 “Communication media” typically embodies computer readable
11 instructions, data structures, program modules, or other data in a modulated data
12 signal, such as carrier wave or other transport mechanism. Communication media
13 also includes any information delivery media.

14 The term “modulated data signal” means a signal that has one or more of its
15 characteristics set or changed in such a manner as to encode information in the
16 signal. By way of example, and not limitation, communication media includes
17 wired media such as a wired network or direct-wired connection, and wireless
18 media such as acoustic, RF, infrared, and other wireless media. Combinations of
19 any of the above are also included within the scope of computer readable media.

20 **Conclusion**

21
22 Although the invention has been described in language specific to structural
23 features and/or methodological steps, it is to be understood that the invention
24 defined in the appended claims is not necessarily limited to the specific features or
25

1 steps described. Rather, the specific features and steps are disclosed as preferred
2 forms of implementing the claimed invention.

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25